## EUROPEAN PATENT APPLICATION

## EUROPEAN PATENT APPLICATION

## INPUT/OUTPUT SYSTEM

This invention relates to an input/output system for single instruction multiple data (SIMD) parallel computers according to the preamble of claim 1 and 14 respectively.

Scientists and engineers from all disciplines have become dependent upon computers to further their work, and with this dependancy they have grown to expect the performance of these computers to increase by an order of magnitude approximately every five years. This trend of increasing computer performance in the order of magnitude range is slowing, in fact, the supercomputers presently available may already be within an order of magnitude of their technological limit. Heretofore, the limit was approximately 3 gigaflops which corresponds to approximately 3 billion floating point instructions per second, which is a function of the length of time it takes electrical signals to propagate through various wires and interconnections at approximately one half the speed of light. The drawback of the prior art system is that many of the problems facing todays scientists and engineers can only be solved utilizing computers with performance capabilities far exceeding the 3 gigaflop limit.

Recent advances in supercomputer performance have been achieved by dividing applications among many processors working in parallel. Theoretically, parallel processing computers should provide performance in the teraflop range. While these computers provide increased capacity and speed, they also provide a new set of problems, namely, programming the new computers, handling the input/output operations and manipulating the data. The programming difficulties stem from the fact that no matter how well a program is written, it is extremely hard to achieve 100 percent utilization of multiple processors. The problem of handling input/output (I/O) operations and data manipulation arises because of the sheer volume of data associated with these types of computers. The programming problem may resolve itself with experience while the I/O and data manipulation problems can be lessened by improving the input/output systems for the computers.

As shown in Fig. 1, a conventional SIMD (single instruction multiple data) parallel system includes a SIMD computer 10 interacts with a host computer 20 via an I/O subsystem 30. The SIMD computer 10 consists of a processor array 11, that includes a plurality of processors 12, numbered P1, P2...PN, each of which is a very simple CPU, a network 13 to connect the processors 12, a memory 14 for each processor, numbered M1, M2...MN, and a control unit 15 to issue instructions and clock

pulses to the processors. The I/O subsystem 30, typically comprises a staging memory that is responsible for transferring data between the SIMD computer 10 and the host 20.

In fine-grained, massively parallel SIMD systems, one single instruction after another is broadcast simultaneously to the processor array, with each instruction being applied to different pieces of data.

Traditionally, fine grained SIMD parallel systems devoted their application emphasis to image-oriented computing which resulted in the input/output system being designed only to handle regularly structured two-dimensional data such as image or matrix data. The input/output rate of a SIMD computer system was typically low due to the fact that for a N-processor SIMD system, arranged as a √N x √N mesh, only √N items of data are input or output to or from the system per machine cycle. Most fine grained SIMD parallel systems are connected by mesh networks and their input/output is done by shifting data between a host and one boundary row/column of the SIMD system. This type of data transfer is considered one dimensional. In addition, data must be pre-arranged by the host such that a particular datum can be assigned to a desired processor. The low input/output rate and restricted capability in handling only regular data structures effectively confine SIMD computers to a narrow application domain.

A second disadvantage of the mesh oriented row/column shifting scheme used in the prior art SIMD input/output systems is the difficulty in programming. Since the input/output function is overlapped with the current task execution, the programmer must interleave the instructions for computing with the instructions for input/output. This situation may lead to a very unreadable code as well as force the programming to stay at the assembly language level.

A third aspect of the prior art input/output subsystems presently employed by SIMD computers is the handling of the corner turning function. The corner turning function is a phenomenon due to the different arrangement of data at the host and SIMD systems. For example, N 32-bit words are arranged in the host as N consecutive words, each being 32-bits wide. However, in transfer, these data words are distributed among 32 planes of SIMD memory with each plane containing N bits, each of which is associated with one processor. This situation arises due to the fact that in the SIMD system, all processors need to access the same memory location in the same machine cycle and the plane organization

supports such memory accessing. The corner-turn-ing of regular data structures such as image or matrix is supported by mesh-oriented row/column shifting. However, corner- turning for irregular data structures is not supported by the prior art row/column shifting I/O scheme.

As noted above, prior art input/output systems are presently implemented as a centralized piece of hardware, such as a staging memory. This approach requires the centralized input/output system to connect to all processors and as a result, many wires are needed for the input/output system. U.S. Patent No. 4,727,474 to Batcher discloses a stag-ing memory for a massively parallel computer. The staging memory is a very complex interface be-tween host memory and local processor memory. This network is capable of buffering, permutating, and shuffling of the data. The circuitry to imple-ment this scheme is complex, requires several stages and is not easily distributed to a very large number of processors.

The mesh-oriented row/column shifting scheme is a compromise, because it connects the input/output system to the boundary of the mesh in order to save wires, but, this in turn, reduces the input/output rate of the system.

U.S. Patent No. 4,380,046 to Fung discloses a massively parallel processor computer which uti-lizes a one-dimensional input/output scheme. The disclosed input/output system serves as a storage element for input and output operations. The in-stantaneous logical state of a bidirectional data bus utilized by the system can be stored into the input/output system in a one bit register and simi-larly, the logical state of the one bit register can be read out to the data bus. The disclosed input/output system is capable of shifting bits to the input/output system in neighboring processing elements. The bits are shifted only in a single direction and thus for a mxm processing element array, one bit slice data stream array will require m shifting operations to move the data array into the processing element array. Thus, there is a need for an I/O system that reduces wiring complexity while maintaining a high input/output rate.

The object of the present invention is to pro-vide an input/output system for a massively parallel SIMD computer with a two-dimensional data trans-fer scheme between a host computer and the SIMD computer, where the SIMD computer is a single instruction, multiple data computer having a parallel array processor comprising a plurality of parallel linked processors each being associated with one of a plurality of SIMD memory devices.

The solution of the objective is described in the characterizing part of claim 1 and 4 for the system and in the characterizing part of claim 14 for the method.

The input/output system comprises a tempo-rary storage means for the bi-directional, two-di-mensional transfer of data between the host com-puter and the SIMD computer, and an input/output processing means for controlling the flow of data between the host computer and the temporary stor-age means, and for controlling the flow of data between the temporary storage means and the plurality of SIMD memory devices. The temporary storage means comprises, in an illustrative embodi-ment of the invention, a plurality of buffers with each one of the plurality of buffers being directly associated with one of the plurality of the SIMD memory devices, and a control circuit means for providing timing and selection signals for the trans-fer of data between the host computer and the temporary storage means and also between the temporary storage means and the SIMD memory devices. The temporary storage means accom-plishes the transfer of data by distributing the data over the plurality of buffers in a predetermined two-dimensional pattern and arranging the data in a format suitable for transfer, in a single system clock cycle.

The input operation of the input/output system of the present invention is a two step process which involves the transfer of data from the host computer memory to the plurality of buffers in the first step and the transfer of data from the plurality of buffers to the SIMD memory devices in the second step. For the transfer of data from the host computer to the plurality of buffers, the input/output processing means writes the I/O data pointer, which is the starting address of the block of data in the host memory to be transferred, and the I/O data length, which is the total number of items to be transferred, to the input/output device of the host computer. Upon completion of the I/O data pointer and the I/O data length transfer, the input/output processing means invokes the transfer of data. The block of data from the host computer memory is distributed to M segments of continuous buffers of the plurality of buffers by having M pairs of segment starting addresses and segment lengths loaded into an address queue of an ad-dress generation unit located in the I/O processing means. The manipulation and control of this data transfer is accomplished by the input/output pro-cessing means and the control circuit means. For the transfer of data from the plurality of buffers to the SIMD memory devices, the input/output pro-cessing means loads the starting address of the SIMD memory devices and length into the address generation unit and then invokes the transfer of data. Once again, the manipulation and control of this data transfer is accomplished by the input/output processing means and the control cir-cuit means. The transfer of data between the host

computer memory and the plurality of buffers is accomplished over the input/output channel while the transfer of data between the plurality of buffers and the plurality of SIMD memory devices is done by a local data bus.

The output operation of the input/output system of the present invention is also a two step process which involves the transfer of data from the plurality SIMD memory devices to the plurality of buffers in the first step and the transfer of data from the plurality of buffers to the host computer memory in the second step. The output operation requires the reverse action and functions of the input operation.

The input/output system of the present invention provides, for a N-processor system, a two-dimensional input/output scheme that supports an input/output rate at a factor of √N higher than the row/column shifting input/output systems utilized in the prior art. The two-dimensionality allows for the efficient transfer of regular data structures as well as for the transfer of irregular data structures, such as sparse matrix or graphic data. This capability permits a user to map data into the processor in an arbitrary predetermined pattern. The present invention is also a distributed architecture which reduces wiring complexity between the input/output system and the SIMD computer. In addition, the input/output system separates input/output programming from computing which reduces the programming effort for a parallel system.

The present invention finds utility in that by incorporating the temporary storage means as an integral and distributed component of the input/output system, two-dimensional data transfer can be accomplished thereby increasing the input/output data rate from √N bits/cycle to N bits/cycle. This type of input/output system greatly increases operating efficiency of any SIMD computer system and can be employed in a plurality of SIMD computer systems since it is independent of the network connecting the processors. The addressing scheme utilized by the input/output system allows the present invention to be utilized by networks using mesh, polymorphic-torus, hypercube and other network connection topologies.

Fig. 1 is a block diagram of a prior art SIMD computer system.

Fig. is a block diagram of a SIMD computer system with one representation of the input/output system of the present invention.

Fig. 3 is a block diagram of a SIMD computer system with another representation of the input/output system of the present invention.

Fig. 4 is a detailed block diagram of a SIMD computer system with another representation of the input/output system of the present invention.

Fig. 5 is a detailed block diagram of the temporary storage means of the present invention.

Fig. 6 is a representation of the mapping scheme for the transfer of data by the input/output system of the present invention.

The input/output system for a massively parallel SIMD computer system is responsible for transferring data between the SIMD computer and its host. Fig. 2 illustrates the basic blocks of a SIMD computer system. The system includes, a host computer 200, which can be a main frame computer or a microprocessor and associated memory, a SIMD computer 100, and an input/output system 300 connecting the host computer 200 and the SIMD computer 100. The input/output system 300 of the present invention provides for the bi-directional, two-dimensional transfer of data between the host computer 200 and the SIMD computer 100.

The SIMD computer 100 comprises a processor array 110 having a plurality of processing elements 120, numbered P1, P2...PN, a network 130 which connects the individual processing elements 120 and a plurality of SIMD memory devices 140, numbered M1, M2...MN. The SIMD computer 100 is a parallel array processor having a great number of individual processing elements 120 linked and operated in parallel. The SIMD computer 100 is massively parallel in that the number N of processing elements 120 is very high, which can be, for example, over one million individual processing elements. The SIMD computer 100 includes a control unit 150 that generates the instruction stream for the processing elements and also provides the necessary timing signals for the computer. The network 130 is an interconnection means for the individual processing elements 120 and can take on many topologies such as mesh, polymorphic-torus and hypercube. The plurality of memory devices 140 are for the immediate storage of data for the individual processing elements 120 and there is a one-to-one correspondence between the number of processing elements 120 and the number of memory devices 140.

The input/output system 300 of the invention includes a temporary storage means 310 coupled to an input/output processor (IOP) 320. The two-dimensional data transfer scheme of the I/O system 300 is provided by the temporary storage means 310. In the illustrative embodiment of Figure 2, the temporary storage means 310 includes a plurality of buffers 330, numbered B1, B2...BN. Each one of the plurality of buffers 330 is associated with one of the plurality of SIMD memory devices 140. The I/O system of the present invention thus utilizes a distributed approach by dividing the I/O data transfer function into N pieces, one for each processor 120. The data to be transferred by the temporary storage means 310 is distributed over said plurality of buffers in a predetermined two-dimensional pat-

tern, and the data is also arranged in a format suitable for transfer, on a single system clock cycle.

The I/O system 300 of the present invention may be configured as a separate entity as in Fig. 2, or the individual elements may be incorporated into other SIMD system components. For example, the IOP functions may be performed by the host 200 and/or the temporary storage means may be incorporated directly in the SIMD processor array 110. Fig. 3 is a block diagram of a SIMD system utilizing both of the above options.

Referring now to Fig. 4, there is shown a detailed diagram of another embodiment of a SIMD system having a host computer 200, a SIMD computer 100 which includes the temporary storage means 310 of the input/output system of the invention incorporated therein and IOP 320 as a separate element. The I/O system further includes an input/output channel 340 which is utilized for the transfer of data between the SIMD computer 100 and the host computer 200. The input/output channel 340 is a n-bit bi-directional data bus which interconnects the host computer 200 and the input/output processing means 320, the host computer 200 and the temporary storage means 310 and the host computer 200 and the array control unit 150. The n-bit bi-directional data bus 340 is capable of handling a multiplicity of data word types depending on the application. For example, the I/O channel 340 may handle single bit, eight bit, sixteen bit and thirty-two bit data words. The input/output processing means 320 controls the overall flow of data in and out of the SIMD computer 100 as well as the flow of data within the computer 100. The input/output processing means 320 is a processor comprising an address generation unit 350, an address queue 360 and a microprocessor and associated memory 370.

The input/output system as stated above is a device capable of the bi-directional two-dimensional transfer of data. Inputting data is accomplished by transferring data from the memory of the host computer 200 to the temporary storage means 310, and then from the temporary storage means 310 to the plurality of SIMD memory devices 140. The outputting of data is accomplished in a similar two-step process wherein the order of the steps comprising the inputting of data is reversed.

INPUTTING DATA FROM HOST TO TEMPORARY STORAGE

To transfer data from the memory of the host computer 200 to the plurality of buffers 330 which comprise the temporary storage means 310, the input/output processor 320 writes the "I/O data

pointer", which is the starting address of the data in the memory of the host computer 200 and the "I/O data length", which is the length of the data in 32-bit words, to the I/O device of the host computer 200 which can be any type of I/O device such as a disk drive or a direct memory access device. Upon completion of this transfer of information, the input/output processor 320 invokes the transfer of data from the memory of the host computer 200 to the temporary storage means 310. The microprocessor and memory 370 contains the I/O program responsible for generating the "I/O data pointer" and the "I/O data length" as well as the necessary instructions for invoking the transfer.

The address generation unit 350 is responsible for generating the address for the particular buffers 330. The input/output processor 320 loads a "segment starting address" and a "segment length" into the address queue 360 of the address generation unit 350 and then invokes the address generation unit 350 and the host input/output device simultaneously for the transfer of data. The address generation unit 350 and the I/O device of the host computer 200 must be synchronized for each datum transfer. The address queue 360 is a first in, first out (FIFO) queue capable of storing multiple segments of addresses. For a continuous block of data in the memory of the host computer 200, the data is distributed to M segments of continuous buffers 330. For this transfer, the input/output processor 320 loads M pairs of "segment starting addresses" (SA) and "segment lengths" (L) into the address queue 360 of the address generation unit 350. The sum of the "segment lengths" is equal to the "I/O data length" written to the host input/output device. In response to receiving the M pairs of "segment starting addresses" and "segment lengths", the address generation unit 350 generates the following addresses:

$$SA(1), SA(1)+1, \ldots, SA(1)+L(1)-1, \qquad (1)$$
$$SA(2), SA(2)+1, \ldots, SA(2)+L(2)-1, \qquad (2)$$
.
.
.
$$SA(M), SA(M)+1, \ldots, SA(M)+L(M)-1. \qquad (3)$$

There are certain possible situations or scenarios where the above described transfer procedure is not straight forward; namely, when the block of data to be transferred has an "I/O data length" larger than the given number of buffers and when the block of data has a word size greater than the buffer width typically 32 bits. To transfer a block of data where the "I/O data length" is larger than the given number of buffers, the input/output processor 320 invokes a program run by microprocessor 370 which transfers the entire data block in several steps. The program ensures that in each step the maximum size of the data transfer is smaller than

the number of buffers 330. To transfer a block of data with word size greater than the buffer width, the host computer 200 must prepare the data so that word size is no greater than 32.

A third situation that arises with the transfer of data is that of having a data word that is smaller than the buffer width. In this case, the data with a word size smaller than the buffer width can be packed into a 32-bit word in the memory of the host computer 200 and distributed to multiple buffers in a single transfer. For example, four bytes of data can be packed into a single 32-bit word and distributed to four continuous buffers in one single transfer. For such a transfer, the input/output processor 320 loads "segment starting address", "segment length", and also "data type" into the address queue 360 of the address generation unit 350. From this input information, the address generation unit 350 generates ADDRESS.BUFFER and ADDRESS.DATATYPE signals which are carried by signal bus 380 to the temporary storage means 310. ADDRESS.BUFFER is a signal which indicates the identifying number of a particular buffer 330 and ADDRESS.DATATYPE is a two bit information code that indicates how many bits are in a particular data word. The code for ADDRESS.DATATYPE may be as follows: 00 denotes that the data being transferred is single bit type, 01 denotes that the data being transferred is eight bit type, 10 denotes that the data being transferred is sixteen bit type, and 11 denotes that the data being transferred is thirty-two bit type. The temporary storage means 310 is responsible for decoding both the AD-DRESS.BUFFER and ADDRESS.DATATYPE. Decoding the ADDRESS.DATATYPE may lead to multiple addressed buffers, for example, in a transfer involving four bytes of data packed into a single 32-bit word, the last two bits of ADDRESS.BUFFER is treated as "don't care", therefore, four buffers are decoded to accept the data. The same 32-bit word is written into four buffers in the same machine cycle. The input/output program executed by microprocessor 370 then rotates the second byte, third byte and the fourth byte into the proper location. The decoding for other data types are performed in a similar manner and the input/output process is completed with the aid of the input/output program contained in microprocessor 370.

Turning now to Fig. 5, there is shown a detailed block diagram of one embodiment of the temporary storage means 310. The temporary storage means 310 is shown consisting of the plurality of buffers 330 and its fundamental support components or circuits as well as the address generation unit 350 which supplies two command signals and the SIMD memory 140. The fundamental components are an address decoder 311, a multiplexing

circuit means 314 which consists of N multiplexers denoted as MUX1 through MUXN, a demultiplexing circuit means 318 consisting of N demultiplexors denoted as DMUX1 through DMUXN, a counter circuit 316 and a comparator circuit 317. Each of the components is fully explained in subsequent paragraphs in conjunction with a description of the operation of the storage means 310. As was stated previously, the address generation unit 350 outputs ADDRESS.BUFFER and ADDRESS.DATATYPE to the temporary storage means 310. These two signals enter the temporary storage means 310 and are decoded by address decoder 311. The address decoder 311 generates a plurality of enable signals, given by

$$EN(i,j).k \qquad (4)$$

where

$$1 \le i \le \sqrt{N}, \qquad (5)$$
$$1 \le j \le \sqrt{N}, \qquad (6)$$

and

$$1 \le k \le 32. \qquad (7)$$

The matrix space defined by i and j represent the total number of buffers and k represents the total capacity of a particular buffer. The total number of buffers in the system is equal to N; therefore, the total number of enable signals is 32xN. Each enable signal represented by equation (4) and carried on line 312 controls the loading of the associated buffer location 330 (B1, B2...BN). When the enable signal is a logic one or in a high state, the associated buffer location is enabled for loading or storing; otherwise, disabled. For AD-DRESS.DATATYPE equal to 11 (32 bit datatype), 32 enable signals, EN(s,t).r are at a high state where s is given by

$$s = ADDRESS.BUFFER/\sqrt{N}, \qquad (8)$$

and t is given by

$$t = ADDRESS.BUFFER-(N*s) \qquad (9)$$

Note that the division by the √N in equation (8) is an integer division which results in the truncation of the remainder of the division.

For ADDRESS.DATATYPE equal to 10 (16 bit datatype), EN(s,t1).r1 and EN(s,t2).r2 are at high states where s is given by

$$s = ADDRESS.BUFFER/\sqrt{N} \qquad (10)$$

t1 is given by

$$t1 = ADDRESS.BUFFER-(s*N), \qquad (11)$$

t2 is given by

$$t2 = t1 + 1, \qquad (12)$$

and r1 and r2 are given by

$$r1 = r2 = 1,2,...16 \qquad (13)$$

For buffer datatype equal to 01 (e.g. byte datatype), four bytes of data will be written into 4 contiguous buffer locations starting from ad-dress.buffer; and for bit datatype (i.e. buffer datatype equal to 00), 32 continuous buffer locations will be selected (neglecting 5 LSB bits of

address.buffer). The calculation of the enable signals are similar to that of the byte datatype.

The address decoder 311 accepts the AD-DRESS.BUFFER and ADDRESS.DATATYPE from the input/output processor 320 and generates the plurality of enable signals. This procedure is used to load the data from the host computer 200 into the buffers 330. Bascially, the data from the host computer 100 is distributed as n-bit words with N addresses. Fig. 6 illustrates the two-dimensional mapping scheme of the present invention. As is shown in this Figure and stated above, the data from the host computer is distributed over the plurality of buffers as n-bit words with N addresses. Each buffer, denoted as B1 through BN represent the starting address for each of the n-bit words. In the illustrative embodiment of the invention, n can be a single bit, eight bits, sixteen bits and thirty-two bits. By generating all the enable signals for a given n-bit word of data, the transfer of the data from the host computer 200 is accomplished in a single system clock cycle. The next step in the process is to transfer the data from the buffers 330 to the SIMD memory devices 140 which occurs in the next system clock cycle.

## INPUTTING DATA FROM TEMPORARY STORAGE TO SIMD MEMORY

Referring once again to Fig. 4, the plurality of SIMD memory devices 140 are shown connected between the temporary storage means 310 and the SIMD processing elements 120. The SIMD memory devices 140 comprises the memory area that interfaces with the buffers 330 and is separately addressable by the address generation unit 350. The SIMD memory is organized as a N-bit wide and D-bit deep memory where N is the total number of processors in the system and D is a choice of implementation. The SIMD memory can be.viewed as D planes each of which consists of N bits of memory. Each bit in a particular plane is denoted as ADDRESS.EXTMEM.BIT which ranges from 0,1...N-1.

For this transfer, the N buffers 330 are organized as 32 planes each containing N bits. Each buffer plane is addressed by AD-DRESS.BUFFER.PLANE. For each system clock cycle, the bits at ADDRESS.BUFFER.PLANE of the buffer at ADDRESS.BUFFER are transferred to the bits at ADDRESS.EXTMEM.BIT of SIMD memory at ADDRESS.EXTMEM. The input/output processor 320 is responsible for inputting the data from the buffers to the SIMD memory. The input/output processor 320 loads the address generation unit 350 with "SIMD memory starting address" and "SIMD length" then invokes the address generation unit 350 to start the transfer.

Referring now to Fig. 5, the exact mechanism for the transfer is explained. A multiplexer/demultiplexer means 314 contains N 32-to-1 multiplexers 315 which selects one out of the 32 locations of the N buffers. All the multiplexers 315 together provide N bits to the plurality of SIMD memory devices 140. The selection control of the multiplexers 315 is provided by a counter means 316 which consists of one 5-bit counter. The 5-bit counter is reset to 0 by the input/output processor 320 upon completion of a write cycle. The counter 316 accepts ADDRESS.DATATYPE from the input/output processor 320 and decodes the ADDRESS.DATATYPE as the length of the word and then stores the length into a comparator 317. For every internal clock cycle, the content of the counter 316 is compared with that of the comparator 317. When equal, a STOP signal is generated to stop the counting, thus indicating that the transfer is complete.

Referring once again to Fig. 6, the n-bit words from the host computer are arranged for transfer to the SIMD memory 140. The first bit location of each buffer, B1 through BN is grouped and denoted as 335(1), the second bit location of each buffer is grouped and denoted as 335(2), and the nth bit location of each buffer is grouped and denoted 335(n). These groupings represent the n planes of memory to be transferred to the SIMD memory 140 from the temporary storage means 330. This Figure represents a grouping of all N buffers; however, as was stated previously, in any particular transfer from the host computer to the temporary storage means, the data is distributed to M segments of buffers where M does not have to correspond to N. Therefore, each of the groupings that represent the n planes of memory need only contain the M locations of data and not N locations. These n planes are addressed by n addresses and each plane contains N bits of data.

## OUTPUTTING DATA

The output operation of the input/output system is also a two-step process, namely, the transfer of data from the SIMD memory 140 to the temporary storage means 310 and the transfer from the temporary storage means 310 to the memory of the host computer 200.

The transfer of data from the SIMD memory to the buffers of the temporary storage means is the reverse action of inputting data from the buffers to the SIMD memory. In the inputting process, n-bit words are written to N addresses by a plurality of multiplexers. In the outputting process N words addressable by n addresses are transferred to the

buffers 330 by means of demultiplexing means 318 which consists of N 1-to-32 demultiplexers 319. The demultiplexers 319 are arranged for transfer to the counters 316 and the comparator 317 in exactly the same manner as described in the inputting process.

The transfer of data from the buffers to the memory of the host computer is the reverse action of inputting from the host to the buffers. In the inputting process, the enable signals determined which buffer can be written to, and in the reverse process, the same enable signals determine which buffers can be read from. The control of this process is by means of the input/output program of the input/output processor.

Turning back to Fig. 6, the n planes of data represented by 335(1) through 335(n) in the SIMD memory 140 are arranged for transfer to the temporary storage means 330. The n planes 335(1) through 335(n) are addressed by N addresses wherein each plane shall contain N addresses for relocation in the temporary storage means.

The concept behind the present invention is a two stage mapping process for the rapid, bi-directional transfer of data between a host computer and a SIMD system. In the transfer of data from the host to the SIMD network the data from the host memory is mapped or distributed over M continuous buffers in a single system clock cycle. In the next clock cycle the data in the M continuous buffers is then distributed over 32 planes of SIMD memory. In the transfer of data from the SIMD network to the host, the data in the SIMD memory is distributed over M continuous buffers in a single system clock cycle. In the next clock cycle the data in M continuous buffers is transferred to the memory of the host computer. As was stated previously, this manipulation of data allows for an increase in data rate of √N for a N processor SIMD system.

The N processors of the SIMD computer can be implemented in a variety of topologies. The preferred topology is to distribute the N processors over a plurality of circuit boards but have collections of processors implemented in a single chip. When each processor in the system is equiped with an associated memory, buffer and a multiplexer/demultiplexer combination and when each collection of processors has an address decoder, a counter and a comparator, then the mapping scheme in Fig. 6 is fully realized. The distributed concept or approach described above has the benefit in VLSI implementation because the wiring between the buffer and the processor/ memory can become intrawire connections within a single chip. This distributed approach greatly reduces the wiring bottleneck in implementing a massively parallel five grained SIMD computer.

Although shown and described in what are believed to be the most practical and preferred embodiments, it is apparent that departures from specific methods and designs described and shown will suggest themselves to those skilled in the art and may be used without departing from the spirit and scope of the invention. The present invention is not restricted to the particular constructions described and illustrated, but should be construed to cohere to all modifications that may fall within the scope of the appended claims.

## Claims

1. Input/output (I/O) system for a massively parallel single instruction multiple data (SIMD) computer providing a two-dimensional data transfer scheme between a host computer and said SIMD computer, said SIMD computer having a parallel array processor comprising a plurality of parallel linked processors each being associated with one of a plurality of SIMD memory devices, characterized by

(a) a temporary storage means (310) coupled between said host computer (200) and said plurality of SIMD memory devices (140) for the bi-directional, two-dimensional transfer of data between said host computer and said SIMD computer;

(b) an input/output processing means (300) for controlling the flow of data between said host computer (200) and said temporary storage means (310), and for controlling the flow of data between said temporary storage means and the plurality of SIMD memory devices (140);

whereby the data to be transferred to and from said temporary storage means is distributed over said temporary storage means in a predetermined two-dimensional pattern, and arranged in a format suitable for transfer, in a single clock cycle.

2. Input/output system of claim 1, characterized in that the temporary storage means (310) includes a plurality of buffers (B1, B2, ...BN), each one of said plurality of buffers being associated with one of said plurality of SIMD memory devices (M1, M2, ...MN).

3. Input/output system of claim 2, characterized in that the temporary storage means (310) includes a control circuit means (311, 350; 314-319) for providing timing and selection signals for the transfer of data between said host computer and said temporary storage means and for the transfer of data between said temporary storage means and said SIMD memory devices.

4. Input/output system according to claim 1 for a massively parallel SIMD computer providing a two-dimensional data transfer scheme between a host computer and said SIMD computer, said SIMD computer having a parallel array processor com-

prising a plurality of parallel linked processors each being associated with one of a plurality of SIMD memory devices, said input/output system comprising:

(a) an input/output channel system (300) for the transfer of data between said SIMD computer (100) and said host computer (200);

(b) a temporary storage (330) connected between said input/output channel and said plurality of SIMD memory devices for the bi-directional, two-dimensional transfer of data between said host computer and said SIMD computer said temporary storage means comprising:

(I) a plurality of buffers (B1-BN), each one of said plurality of buffers being associated with one of said plurality of SIMD memory devices (335(1)-(n), and

(II) a control circuit means (150) for providing timing and selection signals for the transfer of data between said host computer and said temporary storage means and for the transfer of data between said temporary storage means and said SIMD memory devices; and

(c) an input/output processing means (320) for controlling the flow of data between said host computer and said temporary storage means, and for controlling the flow of data between said temporary storage means and said plurality of SIMD memory devices;

whereby the data to be transferred by said temporary storage means is distributed over said plurality of buffers in a predetermined two-dimensional pattern, and arranged in a format suitable for transfer, in a single clock cycle.

5. Input/output system as set forth in claim 4, characterized in
that the input/output channel is a n-bit bi-directional data bus that interconnects said host computer and said input/output processing means, said host computer and said temporary storage means, and said host computer and an array control unit.

6. Input/output system as set forth in claim 5, characterized in
that the temporary storage means is addressable as n-bit words having N addresses, where N equals the number of said buffers and n equals the length of the data words stored in the host memory.

7. Input/output system of claim 6, characterized in
that the control circuit means consists of
multiplexer means (315) for determining which one of n locations of the predetermined number of buffers that data is to be transferred to said plurality of SIMD memory devices;
demultiplexer means (319) for determining which of n locations of the predetermined number of buffers that data is to be transferred into from said plurality of SIMD memory devices;
counter means which provides control signals for

controlling said multiplexer means (316) and said demultiplexer means; and
comparator means (317) for determining the top count for said counter means.

8. Input/output system of claim 7, characterized in
that address decoding means (311) generates said plurality of enable signals from a buffer identification code and a data type code received from said input/output processing means.

9. Input/output system of claim 7, characterized in
that said counter (316) receives said data type code from said input/output processing means and decodes said data type code as the length of the word and stores the length into said comparator means.

10. Input/output system of claim 14, characterized in
that said comparator means (317) compares the count of said counter with said word length and upon a match provides a stop signal to said counter.

11. Input/output system of claim 4, characterized in
that said input/output processing means (320) comprises:
(a) an address generation unit (350) for generating the address for a particular buffer of said plurality of buffers and for generating the address for a particular memory device from said plurality of SIMD memory devices;
(b) a microprocessor and associated memory (370) for generating all control signals for said flow of data; and
an address queue (360) which provides a string of buffer addresses that are sequentially followed.

12. Input/output system as set forth in one of the claim 1-4, characterized by a single instruction multiple data processor (100) comprising:
(a) a parallel array processor (110) comprising a plurality of parallel linked processors (120) each being associated with one of a plurality of SIMD memory devices;
(b) an array control unit (150) for controlling said plurality of parallel linked processors; and
(c) an input/output processor (320) for said single instruction multiple data processor providing a two-dimensional data transfer scheme between a host computer and said array array of arithmetic processing elements, said input/output comprising:
(i) a temporary storage means (370) coupled between said host computer and said SIMD memory devices for the bi-directional, two-dimensional transfer of data between said host computer and said SIMD computer; and
(ii) an input/output processing means (350) for controlling the flow of data between said host computer and said temporary storage

means, and for controlling the flow of data between said temporary storage means and said plurality of SIMD memory devices.

13. Input/output system as set forth in claim 1-4 and 12, characterized in that the single instruction multiple data processor comprises:

(a) a parallel array processor comprising a plurality of parallel linked processors (120) each being associated with one of a plurality of SIMD memory devices;

(b) an array control unit (150) for controlling said plurality of parallel linked processors; and

(c) an input/output system (320) for said single instruction multiple data processor providing a two-dimensional data transfer scheme between a host computer and said array of arithmetic processing elements, said input/output system comprising:

(i) an input/output channel for the transfer of data between said SIMD computer and said host computer;

(ii) a temporary storage means connected between said input/output channel and said plurality of SIMD memory devices for the bi-directional, two-dimensional transfer of data between said host computer and said SIMD computer by arranging and distributing said data over a plurality of buffers in a predetermined pattern in a single system clock cycle, said temporary storage means comprising said plurality of buffers each one of said plurality of buffers being directly associated with one of said plurality of said SIMD memory devices, and a control circuit means for providing timing and selection signals for the transfer of data between said host computer and said temporary storage means and also between said temporary storage means and said SIMD memory devices; and

(iii) an input/output processing means for controlling the flow of data between said host computer and said temporary storage means, and for controlling the flow of data between said temporary storage means and said plurality of SIMD memory devices.

14. Method for the two-dimensional input/output system as set forth in one of claims 1-4, 12 and 13, characterized by

(a) transferring data between a temporary storage means of said SIMD computer and a host computer, said data is transferred utilizing a two-dimensional transfer scheme in a single system clock cycle; and

(b) transferring data between said temporary storage means and a plurality of SIMD memory devices, said data is transferred utilizing said two-dimensional transfer scheme in said single system clock cycle; wherein the step of transfer-

ring data between a temporary storage means and a host computer comprises the steps of:

(a) distributing data from said host computer over a plurality of buffers which comprise said temporary storage means in a first single system clock cycle; and

(b) distributing data from said temporary storage means to a predetermined area of host computer memory in a second single system clock cycle; wherein the step of distributing data from said host computer further includes the steps of:

(a) generating a plurality of enable signals for the transfer of data to a predetermined number of said plurality of buffers; and

(b) transferring the data from said host computer to M segments of contiguous buffers of said plurality of buffers addressable as n-bit words having N addresses, where N equals the number of said plurality of buffers and n is equal to the width of an individual buffer of said plurality of buffers; wherein the step of transferring data between said temporary storage means and a plurality of SIMD memory devices comprises the steps of:

(a) distributing data from said temporary storage means over a plurality of planes which comprise said plurality of SIMD memory devices in a first single system clock cycle; and

(b) distributing data from said plurality of planes over a plurality of buffers which comprise said temporary storage means in a second single system clock cycle.

15. Method according to claim 14, characterized in that the step of distributing data from said temporary storage means further includes the steps of:

(a) determining by means of a multiplexer which one of n locations of said plurality of buffers is to be transferred to said SIMD memory devices; and

(b) transferring the data from said temporary storage means to said plurality of planes addressable as N-bit words having n addresses, where N equals the number of said plurality of buffers and n is equal to the width of an individual buffer of said plurality of buffers; wherein said step of determining the n locations includes the step of generating a set of control signals for controlling said multiplexers.

16. Method according to claim 14, characterized in that the step of distributing data from said plurality of planes further includes the steps of:

(a) determining by means of a demultiplexer which one of n locations of said plurality of buffers data from said plurality of planes is to be

**EP 0 424 618 A2**

transferred into; and

(b) transferring the data from said plurality of planes to M segments of contiguous buffers of said plurality of buffers addressable as n bit words having N addresses, where n equals the number of said plurality of buffers and n is equal to the width of an individual buffer of said plurality of buffers.

FIG. 1

F I G. 2

SIMD NETWORK, 100

ARRAY, 110

M1  M2 --- MN

140

P1  P2 --- PN

120

NETWORK TOPOLOGY

130

ACU —150

I/O SYSTEM, 300

IOP —320

BN  330

TO MN

B2  330

TO M2

B1  330

310

TO M1

200

H O S T

FIG. 3

FIG. 4

INPUT/OUTPUT PROCESSOR

MP

370

360

350  AGU

320

380

ADDR. EXT. MEM.

ADDR. BUFFER

ADDR. DATATYPE

ARRAY CONTROL UNIT

PN

P1

120

MN

M1

140

BN

B1

310

330

ARRAY UNIT

150

110

100

340

340

340

HOST COMPUTER

200

# FIG. 5

FIG. 6



SIMD MEMORY

140

335(n)
335(2)
335(1)

n
2
1

BN
B1

TEMPORARY STORAGE

335(n)
335(2)
335(1)
330

n
2
1

BN
B2
B1

17

Original document

# Input/output system.

| | | |
|---|---|---|
| Publication number: | EP0424618 | Also published as: |
| Publication date: | 1991-05-02 | |
| Inventor: | JAFFE ROBERT S (US); LI HUNGWEN (US); LOHR KIENZLE MARGARET MARY (US); SHENG MING-CHENG (TW) | US5410727 (A |
| | | JP3144783 (A) |
| | | EP0424618 (A |
| Applicant: | IBM (US) | |
| Classification: | | Cited documents: |
| - international: | **G06F15/16; G06F7/78; G06F15/173; G06F15/80; G06F7/76; G06F15/16; G06F15/76;** (IPC1-7): G06F5/06; G06F7/00; G06F13/12; G06F15/16; G06F15/80 | US4727474 |
| | | GB2160685 |
| - european: | | US3287703 |
| Application number: | EP19900115338 19900810 | |
| Priority number(s): | US19890426140 19891024 | |

View INPADOC patent family
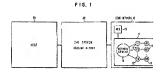
Abstract of **EP0424618**

A two-dimensional input/output system (300) for a massively parallel SIMD computer system (100) providing an interface for the two-way transfer of data between a host computer (200) and the SIMD computer. A plurality of buffers (330) equal in number, and distributed with the individual processing elements (120) of the SIMD computer are used to provide a temporary storage area which allows data in different formats to be mapped in a format suitable for transfer to the host computer or for transfer to the SIMD processing elements. The temporary storage is controlled in such a way as to transfer entire blocks of data in a single SIMD system clock cycle thereby achieving an input/output data rate of N bits/cycle for a SIMD computer consisting of N processors. The system is capable of handling irregular as well as regular data structures. The system also emphasizes a distributed approach in having the input/output system divided into N pieces and distributed to each processor to reduce the wiring complexity while maintaining the I/O rate.

FIG. 1

Data supplied from the *esp@cenet* database - Worldwide

Description of **EP0424618**

INPUT/OUTPUT SYSTEM

This invention relates to an input/output system for single instruction multiple data (SIMD) parallel computers according to the preamble of claim 1 and 14 respectively.

Scientists and engineers from all disciplines have become dependent upon computers to further their wo and with this dependancy they have grown to expect the performance of these computers to increase by order of magnitude approximately every five years. This trend of increasing computer performance in th order of magnitude range is slowing, in fact, the supercomputers presently available may already be wit an order of magnitude of their technological limit. Heretofore, the limit was approximately 3 gigaflops which corresponds to approximately 3 billion floating point instructions per second, which is a function the length of time it takes electrical signals to propagate through various wires and interconnections at approximately one half the speed of light. The drawback of the prior art system is that many of the problems facing todays scientists and engineers can only be solved utilizing computers with performanc capabilities far exceeding the 3 gigaflop limit.

Recent advances in supercomputer performance have been achieved by dividing applications among ma processors working in parallel. Theoretically, parallel processing computers should provide performanc the teraflop range. While these computers provide increased capacity and speed, they also provide a nev set of problems, namely, programming the new computers, handling the input/output operations and manipulating the data. The programming difficulties stem from the fact that no matter how well a progra is written, it is extremely hard to achieve 100 percent utilization of multiple processors. The problem of handling input/output (I/O) operations and data manipulation arises because of the sheer volume of data associated with these types of computers. The programming problem may resolve itself with experience while the I/O and data manipulation problems can be lessened by improving the input/output systems fo the computers.

As shown in Fig. 1, a conventional SIMD (single instruction multiple data) parallel system includes a SIMD computer 10 interacts with a host computer 20 via an I/O subsystem 30. The SIMD computer 10 consists of a processor array 11, that includes a plurality of processors 12, numbered P1, P2...PN, each c which is a very simple CPU, a network 13 to connect the processors 12, a memory 14 for each processo numbered M1, M2...MN, and a control unit 15 to issue instructions and clock pulses to the processors. T I/O subsystem 30, typically comprises a staging memory that is responsible for transferring data betwee the SIMD computer 10 and the host 20.

In fine-grained, massively parallel SIMD systems, one single instruction after another is broadcast simultaneously to the processor array, with each instruction being applied to different pieces of data.

Traditionally, fine grained SIMD parallel systems devoted their application emphasis to image-oriented computing which resulted in the input/output system being designed only to handle regularly structured two-dimensional data such as image or matrix data. The input/output rate of a SIMD computer system v typically low due to the fact that for a N-processor SIMD system, arranged as a 2ROOT N x 2ROOT N mesh, only 2ROOT N items of data are input or output to or from the system per machine cycle. Most fi grained SIMD parallel systems are connected by mesh networks and their input/output is done by shiftir data between a host and one boundary row/column of the SIMD system. This type of data transfer is

considered one dimensional. In addition, data must be pre-arranged by the host such that a particular dat can be assigned to a desired processor.The low input/output rate and restricted capability in handling on regular data structures effectively confine SIMD computers to a narrow application domain.

A second disadvantage of the mesh oriented row/column shifting scheme used in the prior art SIMD input/output systems is the difficulty in programming. Since the input/output function is overlapped witl the current task execution, the programmer must interleave the instructions for computing with the instructions for input/output. This situation may lead to a very unreadable code as well as force the programming to stay at the assembly language level.

A third aspect of the prior art input/output subsystems presently employed by SIMD computers is the handling of the corner turning function. The corner turning function is a phenomenon due to the differer arrangement of data at the host and SIMD systems. For example, N 32-bit words are arranged in the hos as N consecutive words, each being 32-bits wide. However, in transfer, these data words are distributed among 32 planes of SIMD memory with each plane containing N bits, each of which is associated with processor. This situation arises due to the fact that in the SIMD system, all processors need to access the same memory location in the same machine cycle and the plane organization supports such memory accessing. The corner-turning of regular data structures such as image or matrix is supported by mesh-oriented row/column shifting.However, corner- turning for irregular data structures is not supported by prior art row/column shifting I/O scheme.

As noted above, prior art input/output systems are presently implemented as a centralized piece of hardware, such as a staging memory. This approach requires the centralized input/output system to conr to all processors and as a result, many wires are needed for the input/output scheme. U.S. Patent No. 4,727,474 to Batcher discloses a staging memory for a massively parallel computer. The staging memor is a very complex interface between host memory and local processor memory. This network is capable buffering, permutating, and shuffling of the data. The circuitry to implement this scheme is complex, requires several stages and is not easily distributed to a very large number of processors.

The mesh-oriented row/column shifting scheme is a compromise, because it connects the input/output system to the boundary of the mesh in order to save wires, but, this in turn, reduces the input/output rate the system.

U.S. Patent No. 4,380,046 to Fung discloses a massively parallel processor computer which utilizes a or dimensional input/output system. The disclosed input/output system serves as a storage element for inp and output operations. The instantaneous logical state of a bidirectional data bus utilized by the system ( be stored into the input/output system in a one bit register and similarly, the logical state of the one bit register can be read out to the data bus. The disclosed input/output system is capable of shifting bits to tl input/output system in neighboring processing elements. The bits are shifted only in a single direction a thus for a mxm processing element array, one bit slice data stream array will require m shifting operatio to move the data array into the processing element array.Thus, there is a need for an I/O system that reduces wiring complexity while maintaining a high input/output rate.

The object of the present invention is to provide an input/output system for a massively parallel SIMD computer with a two-dimensional data transfer scheme between a host computer and the SIMD compute where the SIMD computer is a single instruction, multiple data computer having a parallel array process comprising a plurality of parallel linked processors each being associated with one of a plurality of SIM memory devices.

The solution of the objective is described in the characterizing part of claim 1 and 4 for the system and i the characterizing part of claim 14 for the method.

The input/output system comprises a temporary storage means for the bi-directional, two-dimensional transfer of data between the host computer and the SIMD computer, and an input/output processing mea for controlling the flow of data between the host computer and the temporary storage means, and for controlling the flow of data between the temporary storage means and the plurality of SIMD memory devices. The temporary storage means comprises, in an illustrative embodiment of the invention, a plurality of buffers with each one of the plurality of buffers being directly associated with one of the plurality of the SIMD memory devices, and a control circuit means for providing timing and selection signals for the transfer of data between the host computer and the temporary storage means and also between the temporary storage means and the SIMD memory devices. The temporary storage means accomplishes the transfer of data by distributing the data over the plurality of buffers in a predetermined two-dimensional pattern and arranging the data in a format suitable for transfer, in a single system clock cycle.

The input operation of the input/output system of the present invention is a two step process which involves the transfer of data from the host computer memory to the plurality of buffers in the first step a the transfer of data from the plurality of buffers to the SIMD memory devices in the second step. For the transfer of data from the host computer to the plurality of buffers, the input/output processing means wr: the I/O data pointer, which is the starting address of the block of data in the host memory to be transferr and the I/O data length, which is the total number of items to be transferred, to the input/output device o the host computer. Upon completion of the I/O data pointer and the I/O data length transfer, the input/output processing means invokes the transfer of data. The block of data from the host computer memory is distributed to M segments of continuous buffers of the plurality of buffers by having M pairs segment starting addresses and segment lengths loaded into an address queue of an address generation u located in the I/O processing means. The manipulation and control of this data transfer is accomplished the input/output processing means and the control circuit means. For the transfer of data from the plural: of buffers to the SIMD memory devices, the input/output processing means loads the starting address of the SIMD memory devices and length into the address generation unit and then invokes the transfer of data. Once again, the manipulation and control of this data transfer is accomplished by the input/output processing means and the control circuit means. The transfer of data between the host computer memory and the plurality of buffers is accomplished over the input/output channel while the transfer of data between the plurality of buffers and the plurality of SIMD memory devices is done by a local data bus.

The output operation of the input/output system of the present invention is also a two step process which involves the transfer of data from the plurality SIMD memory devices to the plurality of buffers in the fi step and the transfer of data from the plurality of buffers to the host computer memory in the second ste] The output operation requires the reverse action and functions of the input operation.

The input/output system of the present invention provides, for a N-processor system, a two-dimensional input/output scheme that supports an input/output rate at a factor of 2ROOT N higher than the row/colur shifting input/output systems utilized in the prior art. The two-dimensionality allows for the efficient transfer of regular data structures as well as for the transfer of irregular data structures, such as sparse matrix or graphic data. This capability permits a user to map data into the processor in an arbitrary predetermined pattern. The present invention is also a distributed architecture which reduces wiring complexity between the input/output system and the SIMD computer. In addition, the input/output syste separates input/output programming from computing which reduces the programming effort for a parall system.

The present invention finds utility in that by incorporating the temporary storage means as an integral ar distributed component of the input/output system, two-dimensional data transfer can be accomplished thereby increasing the input/output data rate from 2ROOT N bits/cycle to N bits/cycle. This type of

input/output system greatly increases operating efficiency of any SIMD computer system and can be employed in a plurality of SIMD computer systems since it is independent of the network connecting the processors. The addressing scheme utilized by the input/output system allows the present invention to be utilized by networks using mesh, polymorphic-torus, hypercube and other network connection topologies.

Fig. 1 is a block diagram of a prior art SIMD computer system.

Fig. is a block diagram of a SIMD computer system with one representation of the input/output system of the present invention.

Fig. 3 is a block diagram of a SIMD computer system with another representation of the input/output system of the present invention.

Fig. 4 is a detailed block diagram of a SIMD computer system with another representation of the input/output system of the present invention.

Fig. 5 is a detailed block diagram of the temporary storage means of the present invention.

Fig. 6 is a representation of the mapping scheme for the transfer of data by the input/output system of the present invention.

The input/output system for a massively parallel SIMD computer system is responsible for transferring data between the SIMD computer and its host. Fig. 2 illustrates the basic blocks of a SIMD computer system. The system includes, a host computer 200, which can be a main frame computer or a microprocessor and associated memory, a SIMD computer 100, and an input/output system 300 connect the host computer 200 and the SIMD computer 100. The input/output system 300 of the present invention provides for the bi-directional, two-dimensional transfer of data between the host computer 200 and the SIMD computer 100.

The SIMD computer 100 comprises a processor array 110 having a plurality of processing elements 120 numbered P1, P2...PN, a network 130 which connects the individual processing elements 120 and a plurality of SIMD memory devices 140, numbered M1, M2...MN. The SIMD computer 100 is a parallel array processor having a great number of individual processing elements 120 linked and operated in parallel. The SIMD computer 100 is massively parallel in that the number N of processing elements 120 very high, which can be, for example, over one million individual processing elements. The SIMD computer 100 includes a control unit 150 that generates the instruction stream for the processing elements and also provides the necessary timing signals for the computer.The network 130 is an interconnection means for the individual processing elements 120 and can take on many topologies such as mesh, polymorphic-torus and hypercube. The plurality of memory devices 140 are for the immediate storage of data for the individual processing elements 120 and there is a one-to-one correspondence between the number of processing elements 120 and the number of memory devices 140.

The input/output system 300 of the invention includes a temporary storage means 310 coupled to an input/output processor (IOP) 320. The two-dimensional data transfer scheme of the I/O system 300 is provided by the temporary storage means 310. In the illustrative embodiment of Figure 2, the temporary storage means 310 includes a plurality of buffers 330, numbered B1, B2...BN. Each one of the plurality buffers 330 is associated with one of the plurality of SIMD memory devices 140. The I/O system of the present invention thus utilizes a distributed approach by dividing the I/O data transfer function into N pieces, one for each processor 120. The data to be transferred by the temporary storage means 310 is distributed over said plurality of buffers in a predetermined two-dimensional pattern, and the data is also arranged in a format suitable for transfer, on a single system clock cycle.

The I/O system 300 of the present invention may be configured as a separate entity as in Fig. 2, or the individual elements may be incorporated into other SIMD system components. For example, the IOP functions may be performed by the host 200 and/or the temporary storage means may be incorporated directly in the SIMD processor array 110. Fig. 3 is a block diagram of a SIMD system utilizing both of the above options.

Referring now to Fig. 4, there is shown a detailed diagram of another embodiment of a SIMD system having a host computer 200, a SIMD computer 100 which includes the temporary storage means 310 of input/output system of the invention incorporated therein and IOP 320 as a separate element. The I/O system further includes an input/output channel 340 which is utilized for the transfer of data between the SIMD computer 100 and the host computer 200. The input/output channel 340 is a n-bit bi-directional d bus which interconnects the host computer 200 and the input/output processing means 320, the host computer 200 and the temporary storage means 310 and the host computer 200 and the array control uni 150. The n-bit bi-directional data bus 340 is capable of handling a multiplicity of data word types depending on the application.For example, the I/O channel 340 may handle single bit, eight bit, sixteen and thirty-two bit data words. The input/output processing means 320 controls the overall flow of data ii and out of the SIMD computer 100 as well as the flow of data within the computer 100. The input/outpu processing means 320 is a processor comprising an address generation unit 350, an address queue 360 a a microprocessor and associated memory 370.

The input/output system as stated above is a device capable of the bi-directional two-dimensional transfe of data. Inputting data is accomplished by transferring data from the memory of the host computer 200 t the temporary storage means 310, and then from the temporary storage means 310 to the plurality of SIM memory devices 140. The outputting of data is accomplished in a similar two-step process wherein the order of the steps comprising the inputting of data is reversed.

INPUTTING DATA FROM HOST TO TEMPORARY STORAGE

To transfer data from the memory of the host computer 200 to the plurality of buffers 330 which compri the temporary storage means 310, the input/output processor 320 writes the "I/O data pointer", which is starting address of the data in the memory of the host computer 200 and the "I/O data length", which is t length of the data in 32-bit words, to the I/O device of the host computer 200 which can be any type of I device such as a disk drive or a direct memory access device. Upon completion of this transfer of information, the input/output processor 320 invokes the transfer of data from the memory of the host computer 200 to the temporary storage means 310. The microprocessor and memory 370 contains the I/ program responsible for generating the "I/O data pointer" and the "I/O data length" as well as the necess instructions for invoking the transfer.

The address generation unit 350 is responsible for generating the address for the particular buffers 330. The input/output processor 320 loads a "segment starting address" and a "segment length" into the addre queue 360 of the address generation unit 350 and then invokes the address generation unit 350 and the h input/output device simultaneously for the transfer of data. The address generation unit 350 and the I/O device of the host computer 200 must be synchronized for each datum transfer. The address queue 360 i first in, first out (FIFO) queue capable of storing multiple segments of addresses. For a continuous block data in the memory of the host computer 200, the data is distributed to M segments of continuous buffer 330.For this transfer, the input/output processor 320 loads M pairs of "segment starting addresses" (SA) and "segment lengths" (L) into the address queue 360 of the address generation unit 350. The sum of the "segment lengths" is equal to the "I/O data length" written to the host input/output device. In response t receiving the M pairs of "segment starting addresses" and "segment lengths", the address generation uni 350 generates the following addresses:
SA(1),SA(1)+1,. . ., SA(1)+L(1)-1, (1)
SA(2),SA(2)+1,. . ., SA(2)+L(2)-1, (2)
.
.

SA(M),SA(M)+1,. . . ., SA(M)+L(M)-1. (3)

There are certain possible situations or scenarios where the above described transfer procedure is not straight forward; namely, when the block of data to be transferred has an "I/O data length" larger than th given number of buffers and when the block of data has a word size greater than the buffer width typica 32 bits. To transfer a block of data where the "I/O data length" is larger than the given number of buffer the input/output processor 320 invokes a program run by microprocessor 370 which transfers the entire data block in several steps. The program ensures that in each step the maximum size of the data transfer smaller than the number of buffers 330. To transfer a block of data with word size greater than the buffe width, the host computer 200 must prepare the data so that word size is no greater than 32.

A third situation that arises with the transfer of data is that of having a data word that is smaller than the buffer width. In this case, the data with a word size smaller than the buffer width can be packed into a 3: bit word in the memory of the host computer 200 and distributed to multiple buffers in a single transfer. For example, four bytes of data can be packed into a single 32-bit word and distributed to four continuo buffers in one single transfer. For such a transfer, the input/output processor 320 loads "segment starting address", "segment length", and also "data type" into the address queue 360 of the address generation ur 350. From this input information, the address generation unit 350 generates ADDRESS.BUFFER and ADDRESS.DATATYPE signals which are carried by signal bus 380 to the temporary storage means 310.ADDRESS.BUFFER is a signal which indicates the identifying number of a particular buffer 330 a ADDRESS.DATATYPE is a two bit information code that indicates how many bits are in a particular d word. The code for ADDRESS.DATATYPE may be as follows: 00 denotes that the data being transferr is single bit type, 01 denotes that the data being transferred is eight bit type, 10 denotes that the data bei transferred is sixteen bit type, and 11 denotes that the data being transferred is thirty-two bit type. The temporary storage means 310 is responsible for decoding both the ADDRESS.BUFFER and ADDRESS.DATATYPE. Decoding the ADDRESS.DATATYPE may lead.to multiple addressed buffer for example, in a transfer involving four bytes of data packed into a single 32-bit word, the last two bits ADDRESS.BUFFER is treated as "don't care", therefore, four buffers are decoded to accept the data.The same 32-bit word is written into four buffers in the same machine cycle. The input/output program executed by microprocessor 370 then rotates the second byte, third byte and the fourth byte into the pro location. The decoding for other data types are performed in a similar manner and the input/output proce is completed with the aid of the input/output program contained in microprocessor 370.

Turning now to Fig. 5, there is shown a detailed block diagram of one embodiment of the temporary storage means 310. The temporary storage means 310 is shown consisting of the plurality of buffers 33( and its fundamental support components or circuits as well as the address generation unit 350 which supplies two command signals and the SIMD memory 140. The fundamental components are an addres: decoder 311, a multiplexing circuit means 314 which consists of N multiplexers denoted as MUX1 thro MUXN, a demultiplexing circuit means 318 consisting of N demultiplexors denoted as DMUX1 throug DMUXN, a counter circuit 316 and a comparator circuit 317. Each of the components is fully explained subsequent paragraphs in conjunction with a description of the operation of the storage means 310.As w stated previously, the address generation unit 350 outputs ADDRESS.BUFFER and ADDRESS.DATATYPE to the temporary storage means 310. These two signals enter the temporary storage means 310 and are decoded by address decoder 311. The address decoder 311 generates a plural of enable signals, given by

EN(i,j).k (4)

where

$1 <\!/= i <\!/= 2\text{ROOT }N$, (5)

$1 <\!/= j <\!/= 2\text{ROOT }N$, (6)

and

$1 </= k </= 32$. (7)

The matrix space defined by i and j represent the total number of buffers and k represents the total capac of a particular buffer. The total number of buffers in the system is equal to N; therefore, the total numbe enable signals is 32xN. Each enable signal represented by equation (4) and carried on line 312 controls t loading of the associated buffer location 330 (B1, B2...BN). When the enable signal is a logic one or in . high state, the associated buffer location is enabled for loading or storing; otherwise, disabled. For ADDRESS.DATATYPE equal to 11 (32 bit datatype), 32 enable signals, EN(s,t).r are at a high state wh s is given by

$s = ADDRESS.BUFFER/ 2ROOT N$, (8)

and t is given by

$t = ADDRESS.BUFFER-(N*s)$ (9)

Note that the division by the 2ROOT N in equation (8) is an integer division which results in the truncat of the remainder of the division.

For ADDRESS.DATATYPE equal to 10 (16 bit datatype), EN(s,t1).r1 and EN(s,t2).r2 are at high states where s is given by

$s = ADDRESS.BUFFER/ 2ROOT N$ (10)

t1 is given by

$t1 = ADDRESS.BUFFER-(s*N)$, (11)

t2 is given by

$t2 = t1 + 1$, (12)

and r1 and r2 are given by

$r1 = r2 = 1,2,...16$ (13)

For buffer datatype equal to 01 (e.g. byte datatype), four bytes of data will be written into 4 contiguous buffer locations starting from address.buffer; and for bit datatype (i.e. buffer datatype equal to 00), 32 continuous buffer locations will be selected (neglecting 5 LSB bits of address.buffer). The calculation o: the enable signals are similar to that of the byte datatype.

The address decoder 311 accepts the ADDRESS.BUFFER and ADDRESS.DATATYPE from the input/output processor 320 and generates the plurality of enable signals. This procedure is used to load t data from the host computer 200 into the buffers 330. Bascially, the data from the host computer 100 is distributed as n-bit words with N addresses. Fig. 6 illustrates the two-dimensional mapping scheme of th present invention. As is shown in this Figure and stated above, the data from the host computer is distributed over the plurality of buffers as n-bit words with N addresses. Each buffer, denoted as B1 through BN represent the starting address for each of the n-bit words. In the illustrative embodiment of t invention, n can be a single bit, eight bits, sixteen bits and thirty-two bits.By generating all the enable signals for a given n-bit word of data, the transfer of the data from the host computer 200 is accomplishe in a single system clock cycle. The next step in the process is to transfer the data from the buffers 330 tc the SIMD memory devices 140 which occurs in the next system clock cycle.

INPUTTING DATA FROM TEMPORARY STORAGE TO SIMD MEMORY

Referring once again to Fig. 4, the plurality of SIMD memory devices 140 are shown connected betwee: the temporary storage means 310 and the SIMD processing elements 120. The SIMD memory devices 1 comprises the memory area that interfaces with the buffers 330 and is separately addressable by the

address generation unit 350. The SIMD memory is organized as a N-bit wide and D-bit deep memory where N is the total number of processors in the system and D is a choice of implementation. The SIMD memory can be viewed as D planes each of which consists of N bits of memory. Each bit in a particular plane is denoted as ADDRESS.EXTMEM.BIT which ranges from 0,1...N-1.

For this transfer, the N buffers 330 are organized as 32 planes each containing N bits. Each buffer plane addressed bY ADDRESS.BUFFER.PLANE. For each system clock cycle, the bits at ADDRESS.BUFFER.PLANE of the buffer at ADDRESS.BUFFER are transferred to the bits at ADDRESS.EXTMEM.BIT of SIMD memory at ADDRESS.EXTMEM. The input/output processor 320 responsible for inputting the data from the buffers to the SIMD memory. The input/output processor 320 loads the address generation unit 350 with "SIMD memory starting address" and "SIMD length" then invokes the address generation unit 350 to start the transfer.

Referring now to Fig. 5, the exact mechanism for the transfer is explained. A multiplexer/demultiplexer means 314 contains N 32-to-1 multiplexers 315 which selects one out of the 32 locations of the N buffer. All the multiplexers 315 together provide N bits to the plurality of SIMD memory devices 140. The selection control of the multiplexers 315 is provided by a counter means 316 which consists of one 5-bit counter. The 5-bit counter is reset to 0 by the input/output processor 320 upon completion of a write cyc. The counter 316 accepts ADDRESS.DATATYPE from the input/output processor 320 and decodes the ADDRESS.DATATYPE as the length of the word and then stores the length into a comparator 317. For every internal clock cycle, the content of the counter 316 is compared with that of the comparator 317.When equal, a STOP signal is generated to stop the counting, thus indicating that the transfer is complete.

Referring once again to Fig. 6, the n-bit words from the host computer are arranged for transfer to the SIMD memory 140. The first bit location of each buffer, B1 through BN is grouped and denoted as 335( the second bit location of each buffer is grouped and denoted as 335(2), and the nth bit location of each buffer is grouped and denoted 335(n). These groupings represent the n planes of memory to be transferr to the SIMD memory 140 from the temporary storage means 330. This Figure represents a grouping of a N buffers; however, as was stated previously, in any particular transfer from the host computer to the temporary storage means, the data is distributed to M segments of buffers where M does not have to correspond to N. Therefore, each of the groupings that represent the n planes of memory need only cont the M locations of data and not N locations.These n planes are addressed by n addresses and each plane contains N bits of data.

OUTPUTTING DATA

The output operation of the input/output system is also a two-step process, namely, the transfer of data from the SIMD memory 140 to the temporary storage means 310 and the transfer from the temporary storage means 310 to the memory of the host computer 200.

The transfer of data from the SIMD memory to the buffers of the temporary storage means is the reverse action of inputting data from the buffers to the SIMD memory. In the inputting process, n-bit words are written to N addresses by a plurality of multiplexers. In the outputting process N words addressable by n addresses are transferred to the buffers 330 by means of demultiplexing 318 which consists of N to-32 demultiplexers 319. The demultiplexers 319 are controlled by the counters 316 and the comparato 317 in exactly the same manner as described in the inputting process.

The transfer of data from the buffers to the memory of the host computer is the reverse action of inputti

from the host to the buffers. In the inputting process, the enable signals determined which buffer can be written to, and in the reverse process, the same enable signals determine which buffers can be read from The control of this process is by means of the input/output program of the input/output processor.

Turning back to Fig. 6, the n planes of data represented by 335(1) through 335(n) in the SIMD memory 140 are arranged for transfer to the temporary storage means 330. The n planes 335(1) through 335(n) a: addressed by N addresses wherein each plane shall contain N addresses for relocation in the temporary storage means.

The concept behind the present invention is a two stage mapping process for the rapid, bi-directional transfer of data between a host computer and a SIMD system. In the transfer of data from the host to the SIMD network the data from the host memory is mapped or distributed over M continuous buffers in a single system clock cycle. In the next clock cycle the data in the M continuous buffers is then distributed over 32 planes of SIMD memory. In the transfer of data from the SIMD network to the host, the data in SIMD memory is distributed over M continuous buffers in a single system clock cycle. In the next clock cycle the data in M continuous buffers is transferred to the memory of the host computer. As was stated previously, this manipulation of data allows for an increase in data rate of 2ROOT N for a N processor SIMD system.

The N processors of the SIMD computer can be implemented in a variety of topologies. The preferred topology is to distribute the N processors over a plurality of circuit boards but have collections of processors implemented in a single chip. When each processor in the system is equiped with an associat memory, buffer and a multiplexer/demultiplexer combination and when each collection of processors ha an address decoder, a counter and a comparator, then the mapping scheme in Fig. 6 is fully realized. The distributed concept or approach described above has the benefit in VLSI implementation because the wiring between the buffer and the processor/ memory can become intrawire connections within a single chip. This distributed approach greatly reduces the wiring bottleneck in implementing a massively paral five grained SIMD computer.

Although shown and described in what are believed to be the most practical and preferred embodiments is apparent that departures from specific methods and designs described and shown will suggest themselves to those skilled in the art and may be used without departing from the spirit and scope of the invention. The present invention is not restricted to the particular constructions described and illustrated but should be constructed to cohere to all modifications that may fall within the scope of the appended claims.

Claims of **EP0424618**

1. Input/output (I/O) system for a massively parallel single instruction multiple data (SIMD) computer providing a two-dimensional data transfer scheme between a host computer and said SIMD computer, s: SIMD computer having a parallel array processor comprising a plurality of parallel linked processors ea being associated with one of a plurality of SIMD memory devices, characterized by
(a) a temporary storage means (310) coupled between said host computer (200) and said plurality of SIM memory devices (140) for the bi-directional, two-dimensional transfer of data between said host comput and said SIMD computer;;
(b) an input/output processing means (300) for controlling the flow of data between said host computer (200) and said temporary storage means (310), and for controlling the flow of data between said tempor

storage means and the plurality of SIMD memory devices (140);
whereby the data to be transferred to and from said temporary storage means is distributed over said temporary storage means in a predetermined two-dimensional pattern, and arranged in a format suitable transfer, in a single clock cycle.

2. Input/output system of claim 1, characterized in
that the temporary storage means (310) includes a plurality of buffers (B1, B2, ...BN), each one of said plurality of buffers being associated with one of said plurality of SIMD memory devices (M1, M2, ...M1

3.Input/output system of claim 2, characterized in
that the temporary storage means (310) includes a control circuit means (311, 350; 314-319) for providi timing and selection signals for the transfer of data between said host computer and said temporary stora means and for the transfer of data between said temporary storage means and said SIMD memory devic

4.Input/output system according to claim 1 for a massively parallel SIMD computer providing a two-dimensional data transfer scheme between a host computer and said SIMD computer, said SIMD compu having a parallel array processor comprising a plurality of parallel linked processors each being associat with one of a plurality of SIMD memory devices, said input/output system comprising:
(a) an input/output channel system (300) for the transfer of data between said SIMD computer (100) anc said host computer (200);
(b) a temporary storage (330) connected between said input/output channel and said plurality of SIMD memory devices for the bi-directional, two-dimensional transfer of data between said host computer and said SIMD computer said temporary storage means comprising::
(i) a plurality of buffers (B1-BN), each one of said plurality of buffers being associated with one of said plurality of SIMD memory devices (335(1)-(n), and
(ii) a control circuit means (150) for providing timing and selection signals for the transfer of data betwe said host computer and said temporary storage means and for the transfer of data between said temporar storage means and said SIMD memory devices; and
(c) an input/output processing means (320) for controlling the flow of data between said host computer a said temporary storage means, and for controlling the flow of data between said temporary storage mear and said plurality of SIMD memory devices;
whereby the data to be transferred by said temporary storage means is distributed over said plurality of buffers in a predetermined two-dimensional pattern, and arranged in a format suitable for transfer, in a single clock cycle.

5. Input/output system as set forth in claim 4, characterized in
that the input/output channel is a n-bit bi-directional data bus that interconnects said host computer and said input/output processing means, said host computer and said temporary storage means, and said hos computer and an array control unit.

6.Input/output system as set forth in claim 5, characterized in
that the temporary storage means is addressable as n-bit words having N addresses, where N equals the number of said buffers and n equals the length of the data words stored in the host memory.

7. Input/output system of claim 6, characterized in
that the control circuit means consists of
multiplexer means (315) for determining which one of n locations of the predetermined number of buff that data is to be transferred to said plurality of SIMD memory devices;
demultiplexer means (319) for determining which of n locations of the predetermined number of buffers that data is to be transferred into from said plurality of SIMD memory devices;
counter means which provides control signals for controlling said multiplexer means (316) and said

demultiplexer means; and
comparator means (317) for determining the top count for said counter means.

8. Input/output system of claim 7, characterized in
that address decoding means (311) generates said plurality of enable signals from a buffer identification
code and a data type code received from said input/output processing means.

9. Input/output system of claim 7, characterized in
that said counter (316) receives said data type code from said input/output processing means and decode
said data type code as the length of the word and stores the length into said comparator means.

10. Input/output system of claim 14, characterized in
that said comparator means (317) compares the count of said counter with said word length and upon a
match provides a stop signal to said counter.

11. Input/output system of claim 4, characterized in
that said input/output processing means (320) comprises:
(a) an address generation unit (350) for generating the address for a particular buffer of said plurality of
buffers and for generating the address for a particular memory device from said plurality of SIMD mem
devices;
(b) a microprocessor and associated memory (370) for generating all control signals for said flow of dat
and
an address queue (360) which provides a string of buffer addresses that are sequentially followed.

12. Input/output system as set forth in one of the claim 1-4, characterized by a single instruction multiple
data processor (100) comprising:
(a) a parallel array processor (110) comprising a plurality of parallel linked processors (120) each being
associated with one of a plurality of SIMD memory devices;
(b) an array control unit (150) for controlling said plurality of parallel linked processors; and
(c) an input/output processor (320) for said single instruction multiple data processor providing a two-
dimensional data transfer scheme between a host computer and said array of arithmetic processing
elements, said input/output comprising:
(i) a temporary storage means (370) coupled between said host computer and said plurality of SIMD
memory devices for the bi-directional, two-dimensional transfer of data between said host computer and
said SIMD computer; and
(ii) an input/output processing means (350) for controlling the flow of data between said host computer
and said temporary storage means, and for controlling the flow of data between said temporary storage
means and said plurality of SIMD memory devices.

13. Input/output system as set forth in claim 1-4 and 12, characterized in that the single instruction multi
data processor comprises:
(a) a parallel array processor comprising a plurality of parallel linked processors (120) each being
associated with one of a plurality of SIMD memory devices;
(b) an array control unit (150) for controlling said plurality of parallel linked processors; and
(c) an input/output system (320) for said single instruction multiple data processor providing a two-
dimensional data transfer scheme between a host computer and said array of arithmetic processing
elements, said input/output system comprising:
(i) an input/output channel for the transfer of data between said SIMD computer and said host computer
(ii) a temporary storage means connected between said input/output channel and said plurality of SIMD
memory devices for the bi-directional, two-dimensional transfer of data between said host computer and
said SIMD computer by arranging and distributing said data over a plurality of buffers in a predetermine

pattern in a single system clock cycle, said temporary storage means comprising said plurality of buffers each one of said plurality of buffers being directly associated with one of said plurality of said SIMD memory devices, and a control circuit means for providing timing and selection signals for the transfer of data between said host computer and said temporary storage means and also between said temporary storage means and said SIMD memory devices; and
(iii) an input/output processing means for controlling the flow of data between said host computer and said temporary storage means, and for controlling the flow of data between said temporary storage means and said plurality of SIMD memory devices.

14. Method for the two-dimensional input/output system as set forth in one of claims 1-4, 12 and 13, characterized by
(a) transferring data between a temporary storage means of said SIMD computer and a host computer, said data is transferred utilizing a two-dimensional transfer scheme in a single system clock cycle; and
(b) transferring data between said temporary storage means and a plurality of SIMD memory devices, said data is transferred utilizing said two-dimensional transfer scheme in said single system clock cycle; wherein the step of transferring data between a temporary storage means and a host computer comprises the steps of:
(a) distributing data from said host computer over a plurality of buffers which comprise said temporary storage means in a first single system clock cycle; and
(b) distributing data from said temporary storage means to a predetermined area of host computer memory in a second single system clock cycle; wherein the step of distributing data from said host computer further includes the steps of:
(a) generating a plurality of enable signals for the transfer of data to a predetermined number of said plurality of buffers; and
(b) transferring the data from said host computer to M segments of contiguous buffers of said plurality of buffers addressable as n-bit words having N addresses, where N equals the number of said plurality of buffers and n is equal to the width of an individual buffer of said plurality of buffers; wherein the step of transferring data between said temporary storage means and a plurality of SIMD memory devices comprises the steps of:
(a) distributing data from said temporary storage means over a plurality of planes which comprise said plurality of SIMD memory devices in a first single system clock cycle; and
(b) distributing data from said plurality of planes over a plurality of buffers which comprise said temporary storage means in a second single system clock cycle.

15.Method according to claim 14, characterized in that the step of distributing data from said temporary storage means further includes the steps of:
(a) determining by means of a multiplexer which one of n locations of said plurality of buffers is to be transferred to said SIMD memory devices; and
(b) transferring the data from said temporary storage means to said plurality of planes addressable as N-words having n addresses, where N equals the number of said plurality of buffers and n is equal to the width of an individual buffer of said plurality of buffers; wherein said step of determining the n location includes the step of generating a set of control signals for controlling said multiplexers.

16. Method according to claim 14, characterized in that the step of distributing data from said plurality of planes further includes the steps of:
(a) determining by means of a demultiplexer which one of n locations of said plurality of buffers data from said plurality of planes is to be transferred into; and
(b) transferring the data from said plurality of planes to M segments of contiguous buffers of said plurality of buffers addressable as n bit words having N addresses, where n equals the number of said plurality of buffers and n is equal to the width of an individual buffer of said plurality of buffers.

Data supplied from the *esp@cenet* database - Worldwide